

AMENDMENTS TO THE CLAIMS

This listing of claims will replace all prior versions, and listings, of claims in the application:

- 1 1. (Cancelled).
- 1 2. (Cancelled).
- 1 3. (Cancelled).
- 1 4. (Cancelled).
- 1 5. (Currently Amended) A computer executable method of managing information about
2 where to begin recovery after a failure of one or more nodes of a multiple-node
3 system, the method comprising the steps of:
4 in a particular node of ~~[[a]]~~ the multiple-node system, maintaining both
5 a single-failure queue ~~that indicates~~ indicating where within a recovery log to
6 begin recovery after a failure of said node; and
7 a multiple-failure queue ~~that indicates~~ indicating where within said recovery
8 log to begin recovery after a failure of said node and one or more other
9 nodes in said multiple-node system;
10 in response to a dirty data item being written to persistent storage, removing an entry
11 for said data item from both said single-failure queue and said multiple-failure
12 queue; and
13 in response to a dirty data item being sent to another node of said multiple-node
14 system without first being written to persistent storage, removing an entry for
15 said data item from said single-failure queue without removing the entry for
16 said data item from said multiple-failure queue.

- 1 6. (Original) The method of Claim 5 further comprising the step of sending the dirty
2 data item to another node to allow removal of the entry from said single-failure queue
3 without the other node requesting the dirty data item.
- 1 7. (Original) The method of Claim 5 further comprising the steps of:
2 after a single node failure, applying said recovery log beginning at a position in said
3 recovery log associated with the single-failure queue; and
4 after a multiple node failure, applying said recovery log beginning at a position in
5 said recovery log associated with the multiple-failure queue.
- 1 8. (Original) The method of Claim 5 wherein:
2 said single-failure queue and said multiple-failure queue are implemented by a single
3 combined queue; and
4 the step of removing an entry for said data item from said single-failure queue
5 without removing the entry for said data item from said multiple-failure queue
6 includes marking an entry for said data item in said combined queue without
7 removing the entry for said data item from said combined queue.
- 1 9. (Original) The method of Claim 5 wherein said single-failure queue and said
2 multiple-failure queue are implemented as two separate queues.
- 1 10. (Currently Amended) A computer executable method for recovering after a failure of
2 one or more nodes of a multiple-node system, the method comprising the steps of:
3 determining whether the failure involves only one node; and
4 if the failure involves only said one node, then performing recovery by applying a
5 recovery log of said node beginning at a first point in the recovery log; and

6 if the failure involves one or more nodes in addition to said one node, then
7 performing recovery by applying said recovery log of said node beginning at a
8 second point in the recovery log;
9 wherein said first point is different from said second point.

1 11. (Original) The method of Claim 10 wherein:

2 the first point is determined, at least in part, by which data items that were dirtied by
3 said node reside in caches in other nodes; and
4 the second point is determined, at least in part, by which data items that were dirtied
5 by said node have been persistently stored.

1 12. (Currently Amended) A computer executable method for recovering after a failure of
2 one or more nodes of a multiple-node system, the method comprising the steps of:

3 if it is unclear whether a particular version of a data item has been written to disk,
4 then performing the steps of
5 without attempting to recover said data item, marking dirtied cached versions
6 of said data item that would have been covered if said particular
7 version was written to disk;
8 when a request is made to write one of said dirtied cached versions to disk,
9 determining which version of said data item is already on disk; and
10 if said particular version of said data item is already on disk, then not writing
11 said one of said dirtied cached versions to disk.

1 13. (Original) The method of Claim 12 further comprising the step of, if said particular
2 version of said data item is not already on disk, then recovering said data item.

1 14. (Original) The method of Claim 12 further comprising the step of, if said particular
2 version of said data item is already on disk, then informing nodes that contain said

3 dirtied cached versions of the data item that said dirtied cached versions are covered
4 by a write-to-disk operation.

1 15. (Cancelled).

1 16. (Currently Amended) A computer executable method for recovering a current
2 version of a data item after a failure of one or more nodes in a system that includes
3 multiple caches, the method comprising the steps of:

4 modifying the data item in a first node of said multiple caches to create a modified
5 data item;

6 sending the modified data item from said first node to a second node of said multiple
7 caches without durably storing the modified data item from said first node to
8 persistent storage;

9 after said modified data item has been sent from said first node to said second node
10 and before said data item in said first node has been covered by a write-to-disk
11 operation, discarding said data item in said first node;

12 after said failure, reconstructing the current version of said data item by applying
13 changes to the data item on persistent storage based on merged redo logs
14 associated with all of said multiple caches;

15 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a
16 locally-dirty checkpoint queue;

17 wherein the globally-dirty data items associated with entries in the globally-dirty
18 checkpoint queue are not retained until covered by write-to-disk operations;

19 determining, for each cache, a checkpoint based on a lower of a first-dirtied time of
20 the entry at the head of the locally-dirty checkpoint queue and the first-dirtied
21 time of the entry at the head of the globally-dirty checkpoint queue; and

after said failure, determining where to begin processing the redo log associated with each cache based on the checkpoint determined for said cache.

17. (Currently Amended) A computer executable method for recovering a current version of a data item after a failure of one or more nodes in a system that includes multiple caches, the method comprising the steps of:

modifying the data item in a first node of said multiple caches to create a modified data item;

sending the modified data item from said first node to a second node of said multiple caches without durably storing the modified data item from said first node to persistent storage;

after said modified data item has been sent from said first node to said second node and before said data item in said first node has been covered by a write-to-disk operation, discarding said data item in said first node; and

after said failure, reconstructing the current version of said data item by applying changes to the data item on persistent storage based on merged redo logs associated with all of said multiple caches;

maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a locally-dirty checkpoint queue;

wherein the globally-dirty data items associated with entries in the globally-dirty checkpoint queue are not retained until covered by write-to-disk operations;

maintaining, for each cache, a first checkpoint record for the locally-dirty checkpoint queue that indicates a first time, where all changes made to data items that are presently dirty in the cache prior to the first time have been recorded on a version of the data item that is on persistent storage;

23 maintaining, for each cache, a second checkpoint record for the globally-dirty
24 checkpoint queue, wherein the second checkpoint record includes a list of data
25 items that were once dirtied in the cache but have since been transferred out
26 and not written to persistent storage; and
27 after said failure, determining where to begin processing the redo log associated with
28 each cache based on the first checkpoint record and said second checkpoint
29 record for said cache.

1 18. (Original) The method of Claim 17 wherein the step of processing the redo log
2 comprises the steps of:
3 determining a starting position for scanning the redo log based on a lesser of a
4 position in the redo log as determined by the first checkpoint record and the
5 positions in the log as determined by the earliest change made to the list of the
6 data items in the second checkpoint record; and
7 during recovery, for the portion of the redo log between the position indicated by the
8 global checkpoint record to the position indicated by the local checkpoint
9 record, considering for potential redo only those log records that correspond to
10 the data items identified in the global checkpoint record.

1 19. (Cancelled).

1 20. (Cancelled).

1 21. (Cancelled).

1 22. (Cancelled).

1 23. (Currently Amended) A computer-readable medium carrying instructions for
2 managing information about where to begin recovery after a failure of one or more

3 nodes of a multiple-node system, the instructions comprising instructions for
4 performing the steps of:
5 in a particular node of a multiple-node system, maintaining both
6 a single-failure queue that indicates where within a recovery log to begin
7 recovery after a failure of said node, and
8 a multiple-failure queue that indicates where within said recovery log to begin
9 recovery after a failure of said node and one or more other nodes in
10 said multiple-node system;
11 in response to a dirty data item being written to persistent storage, removing an entry
12 for said data item from both said single-failure queue and said multiple-failure
13 queue; and
14 in response to a dirty data item being sent to another node of said multiple-node
15 system without first being written to persistent storage, removing an entry for
16 said data item from said single-failure queue without removing the entry for
17 said data item from said multiple-failure queue.

1 24. (Original) The computer-readable medium of Claim 23 further comprising
2 instructions for performing the step of sending the dirty data item to another node to
3 allow removal of the entry from said single-failure queue without the other node
4 requesting the dirty data item.

1 25. (Original) The computer-readable medium of Claim 23 further comprising
2 instructions for performing the steps of:
3 after a single node failure, applying said recovery log beginning at a position in said
4 recovery log associated with the single-failure queue; and

5 after a multiple node failure, applying said recovery log beginning at a position in
6 said recovery log associated with the multiple-failure queue.

1 26. (Original) The computer-readable medium of Claim 23 wherein:

2 said single-failure queue and said multiple-failure queue are implemented by a single
3 combined queue; and

4 the step of removing an entry for said data item from said single-failure queue

5 without removing the entry for said data item from said multiple-failure queue

6 includes marking an entry for said data item in said combined queue without

7 removing the entry for said data item from said combined queue.

1 27. (Original) The computer-readable medium of Claim 23 wherein said single-failure
2 queue and said multiple-failure queue are implemented as two separate queues.

1 28. (Currently Amended) A computer-readable medium carrying instructions for
2 recovering after a failure of one or more nodes of a multiple-node system, the
3 instructions comprising instructions for performing the steps of:

4 determining whether the failure involves only one node; and

5 if the failure involves only said one node, then performing recovery by applying a

6 recovery log of said node beginning at a first point in the recovery log; and

7 if the failure involves one or more nodes in addition to said one node, then

8 performing recovery by applying said recovery log of said node beginning at a

9 second point in the recovery log;

10 wherein said first point is different from said second point.

1 29. (Original) The computer-readable medium of Claim 28 wherein:

2 the first point is determined, at least in part, by which data items that were dirtied by

3 said node reside in caches in other nodes; and

the second point is determined, at least in part, by which data items that were dirtied
by said node have been persistently stored.

30. (Currently Amended) A computer-readable medium carrying instructions for
recovering after a failure of one or more nodes of a multiple-node system, the
instructions comprising instructions for performing the steps of:
if it is unclear whether a particular version of a data item has been written to disk,
then performing the steps of
without attempting to recover said data item, marking dirtied cached versions
of said data item that would have been covered if said particular
version was written to disk;
when a request is made to write one of said dirtied cached versions to disk,
determining which version of said data item is already on disk; and
if said particular version of said data item is already on disk, then not writing
said one of said dirtied cached versions to disk.

31. (Original) The computer-readable medium of Claim 30 further comprising
instructions for performing the step of, if said particular version of said data item is
not already on disk, then recovering said data item.

32. (Original) The computer-readable medium of Claim 30 further comprising
instructions for performing the step of, if said particular version of said data item is
already on disk, then informing nodes ~~that contain~~ containing said dirtied cached
versions of the data item that said dirtied cached versions are covered by a write-to-
disk operation.

33. (Cancelled).

1 34. (Currently Amended) A computer-readable medium carrying instructions for
2 recovering a current version of a data item after a failure of one or more nodes in a
3 system that includes multiple caches, the instructions comprising instructions for
4 performing the steps of:
5 modifying the data item in a first node of said multiple caches to create a modified
6 data item;
7 sending the modified data item from said first node to a second node of said multiple
8 caches without durably storing the modified data item from said first node to
9 persistent storage;
10 after said modified data item has been sent from said first node to said second node
11 and before said data item in said first node has been covered by a write-to-disk
12 operation, discarding said data item in said first node;
13 after said failure, reconstructing the current version of said data item by applying
14 changes to the data item on persistent storage based on merged redo logs
15 associated with all of said multiple caches;
16 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a
17 locally-dirty checkpoint queue;
18 wherein the globally-dirty data items associated with entries in the globally-dirty
19 checkpoint queue are not retained until covered by write-to-disk operations;
20 determining, for each cache, a checkpoint based on a lower of a first-dirtied time of
21 the entry at the head of the locally-dirty checkpoint queue and the first-dirtied
22 time of the entry at the head of the globally-dirty checkpoint queue; and
23 after said failure, determining where to begin processing the redo log associated with
24 each cache based on the checkpoint determined for said cache.

1 35. (Currently Amended) A computer-readable medium carrying instructions for
2 recovering a current version of a data item after a failure of one or more nodes in a
3 system that includes multiple caches, the instructions comprising instructions for
4 performing the steps of:
5 modifying the data item in a first node of said multiple caches to create a modified
6 data item;
7 sending the modified data item from said first node to a second node of said multiple
8 caches without durably storing the modified data item from said first node to
9 persistent storage;
10 after said modified data item has been sent from said first node to said second node
11 and before said data item in said first node has been covered by a write-to-disk
12 operation, discarding said data item in said first node;
13 after said failure, reconstructing the current version of said data item by applying
14 changes to the data item on persistent storage based on merged redo logs
15 associated with all of said multiple caches;
16 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a
17 locally-dirty checkpoint queue;
18 wherein the globally-dirty data items associated with entries in the globally-dirty
19 checkpoint queue are not retained until covered by write-to-disk operations;
20 maintaining, for each cache, a first checkpoint record for the locally-dirty checkpoint
21 queue ~~that indicates~~ indicating a first time, where all changes made to data
22 items that are presently dirty in the cache prior to the first time have been
23 recorded on a version of the data item that is on persistent storage;

24 maintaining, for each cache, a second checkpoint record for the globally-dirty
25 checkpoint queue, wherein the second checkpoint record includes a list of data
26 items that were once dirtied in the cache but have since been transferred out
27 and not written to persistent storage; and
28 after said failure, determining where to begin processing the redo log associated with
29 each cache based on the first checkpoint record and said second checkpoint
30 record for said cache.

1 36. (Original) The computer-readable medium of Claim 35 wherein the step of
2 processing the redo log comprises the steps of:
3 determining a starting position for scanning the redo log based on a lesser of
4 a position in the redo log as determined by the first checkpoint record and
5 the positions in the log as determined by the earliest change made to the list
6 of the data items in the second checkpoint record; and
7 during recovery, for the portion of the redo log between the position indicated by the
8 global checkpoint record to the position indicated by the local checkpoint
9 record, considering for potential redo only those log records that correspond
10 to the data items identified in the global checkpoint record.